

日 本 国 特 許 庁
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出 願 年 月 日

Date of Application:

2002年10月16日

出 願 番 号

Application Number:

特願2002-301539

[ST.10/C]:

[JP 2002-301539]

出 願 人

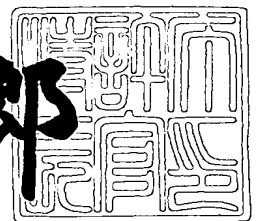
Applicant(s):

インターナショナル・ビジネス・マシーンズ・コーポレーション

2003年 4月11日

特 許 庁 長 官
Commissioner,
Japan Patent Office

太田信一郎



出証番号 出証特2003-3025472

【書類名】 特許願

【整理番号】 JP9020132

【提出日】 平成14年10月16日

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 15/401

【発明者】

 【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ビー・エム株式会社 大和事業所内

 【氏名】 吉田 一星

【特許出願人】

 【識別番号】 390009531

 【氏名又は名称】 インターナショナル・ビジネス・マシーンズ・コーポレーション

【代理人】

 【識別番号】 100086243

 【弁理士】

 【氏名又は名称】 坂口 博

【代理人】

 【識別番号】 100091568

 【弁理士】

 【氏名又は名称】 市位 嘉宏

【代理人】

 【識別番号】 100108501

 【弁理士】

 【氏名又は名称】 上野 剛史

【復代理人】

 【識別番号】 100104880

 【弁理士】

 【氏名又は名称】 古部 次郎

【手数料の表示】

【予納台帳番号】 081504

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9706050

【包括委任状番号】 9704733

【包括委任状番号】 0207860

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 文書自動分類システム、不要語判定方法、文書自動分類方法、およびプログラム

【特許請求の範囲】

【請求項1】 学習用文書集合から語を抽出し、カテゴリごとに語のリストを作成するリスト作成手段と、

前記リスト作成手段により作成された前記リストを用いて、所定の語における各カテゴリでの出現頻度をもとにカテゴリごとの不要語を相対的に決定する不要語決定手段と

を含む文書自動分類システム。

【請求項2】 前記リスト作成手段は、記憶手段にある前記学習用文書集合からカテゴリごとに所定の語における出現頻度を示すリストを生成することを特徴とする請求項1記載の文書自動分類システム。

【請求項3】 前記不要語決定手段は、所定のカテゴリに属する語を取り出し、当該語が他のカテゴリにて所定の基準より多く出現する場合に不要語と決定することを特徴とする請求項1記載の文書自動分類システム。

【請求項4】 前記不要語決定手段は、前記所定のカテゴリから取り出される前記語が、予め定められる閾値および前記他のカテゴリに属する文書の個数によって決定される前記所定の基準より当該他のカテゴリにて多く出現する場合に不要語と決定することを特徴とする請求項3記載の文書自動分類システム。

【請求項5】 前記不要語決定手段により決定され、当該決定により不要語が除去された、カテゴリごとのリストを分類用カタログとして格納する分類用カタログ格納手段と、

前記分類用カタログ格納手段に格納された前記分類用カタログを用いて、分類対象文書に対して分類処理を施す文書分類手段と
を更に含む請求項1記載の文書自動分類システム。

【請求項6】 分野ごとに分類済みの文書を格納する分類済み文書集合格納装置と、

前記分類済み文書集合格納装置から取得された文書に含まれる単語の出現頻度

の情報を含む分野別のテーブルを作成する分野テーブル作成部と、

前記分野テーブル作成部により作成された前記分野別のテーブルから得られる所定の語における各分野での出現頻度に基づいて、当該テーブルから当該分野別に不要語を除去する不要語除去部と、

前記不要語除去部により不要語が除去された前記テーブルを格納する分類用カタログ格納装置と

を含む文書自動分類システム。

【請求項 7】 分類される分類対象文書を格納する分類対象文書格納装置と

前記分類対象文書格納装置に格納された前記分類対象文書に対し、前記分類用カタログ格納装置に格納された前記テーブルを用いて分類処理を行なう文書分類処理部と

を更に含む請求項 6 記載の文書自動分類システム。

【請求項 8】 前記不要語除去部は、所定の分野に属する語を取り出し、当該語が他の分野にて所定の基準を超えて出現する場合に、当該語を不要語として前記テーブルから除去することを特徴とする請求項 6 記載の文書自動分類システム。

【請求項 9】 前記分野テーブル作成部により作成される前記分野別のテーブルは、前記単語、当該単語の出現頻度、および当該単語の品詞に関する情報を含むことを特徴とする請求項 6 記載の文書自動分類システム。

【請求項 10】 文書自動分類システムにおける不要語判定方法であって、学習用文書集合が格納されている記憶装置から文書に含まれる単語をカテゴリごとに抽出するステップと、

抽出された前記単語の出現頻度の情報を含むリストをカテゴリごとに作成するステップと、

作成された前記リストを用いて、所定のカテゴリに属する所定の単語における他のカテゴリでの出現頻度を認識するステップと、

認識された前記出現頻度に基づいて、カテゴリごとに不要語を判定するステップと

を含む不要語判定方法。

【請求項 1 1】 前記不要語を判定するステップでは、前記所定のカテゴリから選定された 1 つの単語に対して、前記他のカテゴリにて当該単語が所定の基準を超えて含まれているか否かによって不要語を判定することを特徴とする請求項 1 0 記載の不要語判定方法。

【請求項 1 2】 前記所定の基準は、前記他のカテゴリ内の文書数および予め定められた所定の閾値により得られる値であることを特徴とする請求項 1 1 記載の不要語判定方法。

【請求項 1 3】 前記所定の基準は、前記他のカテゴリ内での前記単語の頻度と、当該他のカテゴリ内での全ての単語の頻度合計とによって決定されることを特徴とする請求項 1 1 記載の不要語判定方法。

【請求項 1 4】 文書自動分類システムにおける不要語判定方法であって、記憶装置に格納されている分野別に分類済みの文書集合から分野別の単語に関する情報を取得し、

取得された前記単語に関する情報に基づいて、特定の分野に属する単語が他の分野に出現する頻度を認識し、

認識される前記頻度に基づいて前記単語が前記特定の分野を識別するのに不要な単語か否かを判定することを特徴とする文書自動分類方法。

【請求項 1 5】 前記不要な単語であると判定された単語を除去して文書の分類用カタログを生成し、

生成された前記分類用カタログを記憶装置に格納することを特徴とする請求項 1 4 記載の文書自動分類方法。

【請求項 1 6】 前記記憶装置に格納された前記分類用カタログを用いて、分類対象文書に対して分類処理を施すことを特徴とする請求項 1 5 記載の文書自動分類方法。

【請求項 1 7】 コンピュータに、
学習用文書集合が格納されている記憶装置から文書に含まれる単語をカテゴリごとに抽出する機能と、

抽出された前記単語の出現頻度の情報を含むリストをカテゴリごとに作成する

機能と、

作成された前記リストを用いて、所定のカテゴリに属する所定の単語における他のカテゴリでの出現頻度を認識する機能と、

認識された前記出現頻度に基づいて、カテゴリごとに不要語を判定する機能とを実現させるプログラム。

【請求項 1 8】 前記コンピュータに、

判定された前記不要語を用いて分類用のリストを生成する機能を更に実現させる請求項 1 7 記載のプログラム。

【請求項 1 9】 コンピュータに、

記憶装置に格納されている分野別に分類済みの文書集合から分野別の単語に関する情報を取得する機能と、

取得された前記単語に関する情報に基づいて、特定の分野に属する単語が他の分野に出現する頻度を認識する機能と、

認識される前記頻度に基づいて前記単語が前記特定の分野を識別するのに不要な単語か否かを判定する機能とを実現させるプログラム。

【請求項 2 0】 前記コンピュータに、

前記不要な単語であると判定された単語を除去して文書の分類用カタログを生成する機能と、

生成された前記分類用カタログを用いて、分類対象文書を分類する機能とを更に実現させる請求項 1 9 記載のプログラム。

【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術分野】

本発明は、文書データを自動的に分類するための文書自動分類システム等に係り、より詳しくは、不要語を効果的に取り除く文書自動分類システム等に関する。

【0 0 0 2】

【従来の技術】

近年、電子化された文書データ(テキスト)が大量に流通するのに伴い、例えば文書格納データベースに存在する大量の文書を自動分類する文書自動分類システムが注目されている。この文書自動分類システムは、学習機能と分類機能との2つの要素から構成される。これらの機能を実現するために、決定木、Neural Network、ベクトル空間モデルなど、様々なモデルが提案されている。何れの方法においても、各カテゴリや文書を特定付ける語を文書から抽出することが重要である。しかしながら、文書から単語を頻度順に取り出す場合には、カテゴリを一意に決定するために、有用でない語(不要語)が上位を占めてしまう。この不要語を学習・分類前に除去しておくことで、文書自動分類システムの分類性能を大きく改善することができる。

【0003】

不要語には、大きく分けて、機能語と一般語の2種類が存在する。機能語は、語と語の関係を表す助詞、助動詞などを指す。この機能語は、カテゴリに存在しないものが多いので、語の品詞を調べたり、予め不要語リストを作成しておくことにより除去することができる。一方、一般語は、機能語以外に一般的に用いられる語を指す。この一般語は、機能語と異なり、語の頻度によって決定される場合が多く、与えられた文書集合中の出現頻度がある上限または下限を越えた語を不要語とする手法が一般的に用いられる。この上限、下限を決める手法として、語の出現頻度に関する経験法則をもとに、多過ぎる或いは少な過ぎる語を判定して除去するZipfの法則などが知られている。

【0004】

文書自動分類技術に関する従来技術として、例えば、分類済み文書から複数分野語を学習し、この複数分野語に注目して関連度テーブルや分類対象文書の単語の頻度情報を詳細化することで、分類対象文書の各分野への関連度をより詳細に分析し、類似する分野での分類精度を向上させるものが存在する(例えば、特許文献1参照)。また、不要単語を登録した不要語辞書を設け、新規単語に不要語辞書中の不要単語と同一のものが含まれているときに新規単語を削除し、不要単語が削除された新規単語に対して単語重要度を決定する技術について開示されている(例えば、特許文献2参照)。更に、精度の高い類似文書検索を行なうために

、出現頻度をカウントして不要語リストを自動的に作成し、一定の割合以上(以下)出現する語を削除することで、類似度算出精度の向上を図った技術が示されている(例えば、特許文献3参照)。

【0005】

【特許文献1】

特開平10-254883号公報(第4-5頁、15頁、図1)

【特許文献2】

特開平11-120183号公報(第3-4頁、図1)

【特許文献3】

特開平11-259515号公報(第3-5頁、図3)

【0006】

【発明が解決しようとする課題】

このように、精度の高い文書自動分類を実行するためには、文書中に存在する抽出すべき単語対象から不要語を排除することが好ましい。しかしながら、特許文献1では、まず、不要語除去という概念が存在せず、どの単語も最低1つは関連の強いカテゴリが存在することを前提としており、品詞の限定や不要語リストの作成を行わない限り不要語がそのままリスト登録されてしまい、精度の高い分類が困難となる。また、関連テーブルを作成した後、改めて詳細関連度テーブルを作成しており、多くの記憶容量を必要としてしまう。

【0007】

また、特許文献2では、予め用意された不要語リストとの照合による不要語除去がなされているが、対象となるカテゴリの集合ごとに不要語リストを作り直す必要があり、また、時代と共に変化する用語に対して十分に対処することができない。更に、特許文献3では、学習用文書全体における各語の出現頻度をカウントしているものの、頻度の基準値を設定してそれを越えた語を除去する方法に留まり、除去されない不要語が数多く残る可能性が高く、その一方で、不要語の判定を広く行なうと、分類のために有用な語まで除去されてしまうという問題があった。また、上述したZipfの法則では、上限・下限を越えない語の中にも不要語が含まれていたり、逆に、上限・下限を越えた語の中にカテゴリを特徴付け

る重要な語が含まれている場合がある。

【0008】

本発明は、以上のような技術的課題を解決するためになされたものであって、その目的とするところは、文書自動分類において、不要語を効果的に取り除くことにある。

【0009】

【課題を解決するための手段】

かかる目的のもと、本発明は、文書を自動的にカテゴリに分類する文書自動分類システムにおいて、学習用文書集合から語を抽出し、カテゴリごとに語のリストを作成するリスト作成手段と、このリスト作成手段により作成されたリストを用いて、所定の語における各カテゴリでの出現頻度をもとにカテゴリごとの不要語を相対的に決定する不要語決定手段と、この不要語決定手段により決定され、不要語が除去された、カテゴリごとのリストを分類用カタログとして格納する分類用カタログ格納手段と、この分類用カタログ格納手段に格納された分類用カタログを用いて、分類対象文書に対して分類処理を施す文書分類手段とを含む。

【0010】

ここで、このリスト作成手段は、記憶手段にある学習用文書集合からカテゴリごとに所定の語における出現頻度を示すリストを生成することを特徴とすることができる。また、この不要語決定手段は、所定のカテゴリに属する語を取り出し、語が他のカテゴリにて所定の基準より多く出現する場合に不要語と決定すれば、カテゴリ間の相対的な出現頻度をもとに不要語を判断することができ、効果的に不要語を除去することができる。更に、この不要語決定手段は、所定のカテゴリから取り出される語が、予め定められる閾値および他のカテゴリに属する文書の個数によって決定される所定の基準より他のカテゴリにて多く出現する場合に不要語と決定することを特徴とすることができる。

【0011】

他の観点から捉えると、本発明が適用される文書自動分類システムは、分野ごとに分類済みの文書を格納する分類済み文書集合格納装置と、この分類済み文書集合格納装置から取得された文書に含まれる単語の出現頻度の情報を含む分野別

のテーブルを作成する分野テーブル作成部と、この分野テーブル作成部により作成された分野別のテーブルから得られる所定の語における各分野での出現頻度に基づいて、テーブルから分野別に不要語を除去する不要語除去部と、この不要語除去部により不要語が除去されたテーブルを格納する分類用カタログ格納装置と、分類される分類対象文書を格納する分類対象文書格納装置と、この分類対象文書格納装置に格納された分類対象文書に対し、分類用カタログ格納装置に格納されたテーブルを用いて分類処理を行なう文書分類処理部とを含む。

【 0 0 1 2 】

一方、本発明は、文書自動分類システムにおける不要語判定方法であって、分野テーブル作成手段により、学習用文書集合が格納されている記憶装置から文書に含まれる単語をカテゴリごとに抽出し、抽出された単語の出現頻度の情報を含むリストをカテゴリごとに作成するステップと、不要語判定手段により、作成されたリストを用いて、所定のカテゴリに属する所定の単語における他のカテゴリでの出現頻度を認識し、認識された出現頻度に基づいて、カテゴリごとに不要語を判定するステップとを含む。

【 0 0 1 3 】

ここで、この不要語を判定するステップでは、所定のカテゴリから選定された1つの単語に対して、他のカテゴリにてこの単語が所定の基準を超えて含まれているか否かによって不要語を判定することを特徴とすれば、カテゴリの特徴付けに役立たない語を有効に取り除くことができる点で好ましい。また、この所定の基準は、他のカテゴリ内の文書数および予め定められた所定の閾値により得られる値であることを特徴とすることができる。また他の態様として、この所定の基準は、他のカテゴリ内での単語の頻度と、この他のカテゴリ内での全ての単語の頻度合計とによって決定されることを特徴とすることもできる。

【 0 0 1 4 】

更に他の観点から捉えると、本発明が適用される文書自動分類方法は、記憶装置に格納されている分野別に分類済みの文書集合から分野別の単語に関する情報を取得し、取得された単語に関する情報に基づいて、特定の分野に属する単語が他の分野に出現する頻度を認識し、認識される頻度に基づいてこの単語が特定の

分野を識別するのに不要な単語か否かを判定し、不要な単語であると判定された単語を除去して文書の分類用カタログを生成し、生成された分類用カタログを記憶装置に格納し、この記憶装置に格納された分類用カタログを用いて、分類対象文書に対して分類処理を施すことを特徴とすることができる。

【 0 0 1 5 】

尚、これらの発明は、コンピュータに各機能を実現させるプログラムとして機能させることができる。より具体的には、コンピュータに、学習用文書集合が格納されている記憶装置から文書に含まれる単語をカテゴリごとに抽出する機能と、抽出された単語の出現頻度の情報を含むリストをカテゴリごとに作成する機能と、作成されたリストを用いて、所定のカテゴリに属する所定の単語における他のカテゴリでの出現頻度を認識する機能と、認識された出現頻度に基づいて、カテゴリごとに不要語を判定する機能と、判定された不要語を用いて分類用のリストを生成する機能を実現させるプログラムとして把握することができる。

【 0 0 1 6 】

また、本発明は、コンピュータに、記憶装置に格納されている分野別に分類済みの文書集合から分野別の単語に関する情報を取得する機能と、取得された単語に関する情報に基づいて、特定の分野に属する単語が他の分野に出現する頻度を認識する機能と、認識される頻度に基づいて単語が特定の分野を識別するのに不要な単語か否かを判定する機能と、不要な単語であると判定された単語を除去して文書の分類用カタログを生成する機能と、生成された分類用カタログを用いて、分類対象文書を分類する機能とを実現させるプログラムとして把握することができる。

【 0 0 1 7 】

尚、これらのプログラムとしては、コンピュータを顧客に対して提供する際に、コンピュータ装置の中にインストールされた状態にて提供される場合の他、コンピュータに実行させるプログラムをコンピュータが読取可能に記憶した記憶媒体にて提供する形態が考えられる。この記憶媒体としては、例えばCD-ROM媒体等が該当し、CD-ROM読取装置等によってプログラムが読み取られ、フラッシュROM等にこのプログラムが格納されて実行される。また、これらのプ

プログラムは、例えば、プログラム伝送装置によってネットワークを介して提供される形態がある。このプログラム伝送装置としては、例えば、ネット上のサーバに設けられ、プログラムを格納するメモリと、ネットワークを介してプログラムを提供するプログラム伝送手段とを備えている。

【 0 0 1 8 】

【発明の実施の形態】

以下、添付図面を参照し、本発明が適用される実施の形態について詳細に説明する。

図 1 は、本実施の形態が適用される文書自動分類システム 1 0 の構成を示したブロック図である。この文書自動分類システム 1 0 は、パーソナルコンピュータ (P C) 等のコンピュータ装置によって展開され、 H D D (ハードディスクドライブ) などの外部記憶装置にて構成されて各種データを格納するデータ格納装置 2 0 と、外部メモリから読み出されたアプリケーションプログラムによって C P U により実行される処理部 3 0 とを備えている。実際には、処理部 3 0 の各ブロック構成要素は、 C P U の実行プログラムの読み込み領域として或いは実行プログラムの処理データを書き込む作業領域として利用される複数の D R A M チップ等からなる内部記憶装置にて展開される。

【 0 0 1 9 】

データ格納装置 2 0 は、分野 (カテゴリ) の学習処理に用いられる分類済みの文書を格納する、学習用文書集合である分類済み文書集合格納装置 2 1、不要語が除去された後の分類用カタログが格納される分類用カタログ格納装置 2 2、実際に文書分類処理がなされる対象となるテキストを格納する分類対象文書格納装置 2 3、分類された結果を格納する分類結果格納装置 2 4 を備えている。この分類結果格納装置 2 4 の内容は、分類済み文書集合格納装置 2 1 に格納されて学習処理に利用されるように構成することもできる。ここで、「不要語」とは、例えば、カテゴリ (分野) の特徴付けに役立たない語と定義できる。

【 0 0 2 0 】

処理部 3 0 は、不要語削除前に選択された分野 (カテゴリ) ごとに、語のリストであるテーブル情報を生成する分野テーブル作成部 3 1、分野テーブル作成部 3

1にて作成された分野テーブルの単語について、不要語の決定と決定された不要語を除去する処理とを実行する不要語決定・除去部32、実際に文書分類処理を実行する文書分類処理部33を備えている。

【0021】

分野テーブル作成部31は、分類済み文書集合格納装置21から得られた文書を用いて、例えば単語の出現頻度等の情報を含むテーブルを作成し、テーブル情報として内部記憶装置に登録している。分類済み文書集合格納装置21では、学習用文書である複数の文書が、例えば、「政治」「経済」「スポーツ」等の分野(カテゴリ)の集合に分類された状態にて格納されている。分野テーブル作成部31では、このカテゴリの集合に分類された文書を読み込み、その文書を解析し、例えば、文書に含まれる単語(語)の出現頻度をカウントして、分野テーブルを生成している。テーブルのデータ量が多い場合には、外部記憶装置であるデータ格納装置20に別途、格納するように構成することができる。尚、分類済み文書集合格納装置21の代わりに、所定のネットワークを介して、学習用文書集合(分類済み文書集合)を取得するように構成することも可能である。

【0022】

不要語決定・除去部32では、分野テーブル作成部31にて作成された分野テーブルを用いて、カテゴリ間の相対的な出現頻度によって不要語を決定する処理が実行される。不要語決定・除去部32によって不要語が除去された分野テーブルは、分類用カタログ格納装置22に格納される。

【0023】

文書分類処理部33では、分類対象文書格納装置23に格納されている実際の分類対象となる文書に対し、分類用カタログ格納装置22に格納された分類用カタログ(不要語が除去された分野テーブル)を用いて、文書分類処理が実行される。この文書分類処理部33による分類結果は、分類結果格納装置24に格納される。

【0024】

ここで、分野テーブル作成処理について説明する。

図2は、分野テーブル作成部31にてなされる処理を示したフローチャートで

ある。分野テーブルの作成に際し、分野テーブル作成部 3 1 では、分類済み文書集合格納装置 2 1 に格納されている全ての分野について作業がなされているか否かが判断される(ステップ 1 0 1)。全ての分野についての作業がなされていない場合には、まず、分野を 1 つ選び(ステップ 1 0 2)、作業していない文書が分野内にあるか否かの判断がなされる(ステップ 1 0 3)。分野内にない場合には、ステップ 1 0 1 に戻り、まだ残っている場合には、その分野から文書を 1 つ選ぶ(ステップ 1 0 4)。そして、作業していない単語が文書内にあるか否かが判断され(ステップ 1 0 5)、もう残っていない場合には、ステップ 1 0 3 へ戻り、未処理の単語がまだ文書内にある場合には、文書から単語が 1 つ選ばれる(ステップ 1 0 6)。この単語の抽出では、形態素解析が用いられる。また、品詞によるフィルタリングをこのタイミングで行なうこともできる。

【 0 0 2 5 】

そして、単語が既にテーブル(分野テーブル)に登録されているか否かが判断され(ステップ 1 0 7)、登録されている場合には、テーブル上の、登録単語の頻度(出現頻度)を 1 増やして、ステップ 1 0 5 に戻る。登録されていない単語である場合には、その単語をテーブルに登録し(ステップ 1 0 9)、ステップ 1 0 5 へ戻る。このテーブル(分野テーブル)には、単語とその出現頻度の他に、各単語ごとに情報を持つこともできる。例えば、単語の品詞情報などを持つことができ、かかる場合には、この品詞情報等もテーブルに登録される。これらの一連の処理を行い、ステップ 1 0 1 にて全ての分野について作業をしたと判断された場合には、分野テーブル作成処理が終了する。

【 0 0 2 6 】

図 3 は、図 2 にて説明したような分野テーブル作成部 3 1 にて作成されメモリに格納されるテーブル例を示した図である。ここでは、「スポーツ」の分野について、不要語除去前のテーブル例が示されている。テーブル情報は、単語を特定する番号である単語 ID ごとに、単語、単語の品詞、単語の出現頻度が表されている。この単語の出現頻度は、「学習用の文書集合中に出現した回数の総計」を表している。1 文書内に 2 個以上出現した場合も、その個数の分だけ数えている。尚、図 3 に示す例では、予め「名詞」と「動詞」だけをテーブルに登録すると

いう、前処理を行なってできたテーブルの模式図である。

【 0 0 2 7 】

次に、不要語除去処理について説明する。

図 4 は、不要語決定・除去部 3 2 にてなされる処理を示したフローチャートである。不要語決定・除去部 3 2 では、分野テーブル作成部 3 1 にて作成された分野テーブルを用いて、全ての分野について作業がなされているか否かが判断される(ステップ 2 0 1)。全ての分野についての作業がなされていない場合には、まず、1つの分野(Aとする)を選ぶ(ステップ 2 0 2)。そして、Aの分野テーブルにおける全ての単語について作業したか否かが判断され(ステップ 2 0 3)、全ての単語について作業した場合には、ステップ 2 0 1 に戻り、まだ残っている場合には、Aの分野テーブルから1つの単語(Wとする)を選ぶ(ステップ 2 0 4)。そして、A以外の全ての分野と比較したか否かが判断され(ステップ 2 0 5)、A以外の全ての分野と比較した場合には、ステップ 2 0 3 へ戻り、比較していない場合には、A以外の分野から1つの分野(Bとする)を選ぶ(ステップ 2 0 6)。そして、予め定められている基準に対し、この基準を超えて、Bの分野テーブルにWが含まれているか否かが判断され(ステップ 2 0 7)、基準を超えて含まれていない場合には、ステップ 2 0 5 からの処理に戻る。含まれている場合には、このWを不要語と判定して(ステップ 2 0 8)、ステップ 2 0 3 からの処理に戻る。ステップ 2 0 1 にて、全ての分野について作業がなされたと判断された場合には、不要語除去の処理は終了し、除去結果のテーブル情報が分類用カタログ格納装置 2 2 に格納される。

【 0 0 2 8 】

即ち、図 4 に示す不要語除去の方法では、所定のカテゴリ A に属する語(単語) W を 1 つ、取り出し、この語 W が他のカテゴリ B の中で、所定の基準より多く出現するならば、この語 W がカテゴリ A の不要語と決定している。これを、カテゴリ A に属する全ての語について行なう。また、この一連の処理を、カテゴリ A 以外のカテゴリに対して、それぞれの判断対象となるカテゴリの役割を取り替えて、全てのカテゴリに対して不要語を決定している。

【 0 0 2 9 】

ここで、ステップ207の判断である「基準を超えて含まれている」を定義する方法としては、幾つかの方法が考えられる。例えば、後述するように、ある閾値を定め、分類済み文書集合格納装置21に格納された学習用文書の個数に対し

文書数×閾値

で得られた値に対してB内での単語Wの頻度が超える場合には、「基準値を超えて含まれる」と定義することができる。また、他の例として、例えば、

単語WのB内での頻度 ÷ B内の全ての単語の頻度合計

がある閾値を超えた場合には、「基準を超えて含まれる」と定義するように構成することもできる。

【0030】

尚、図4に示す不要語除去の方法は、他の既存の不要語除去手法と組み合わせて用いることも可能である。また、分野(カテゴリ)が階層構造を成している場合も、同一階層に存在する分野に対してこのアルゴリズムを適用することによって、拡張することもできる。

【0031】

図5(a)～(c)は、この不要語処理のアルゴリズムについて更に詳しく説明するための図である。本アルゴリズムでは、まず、閾値 $R(0 \leq R \leq 1)$ が処理部30に格納される。図5(a)～(c)に示す例では、この閾値として、「0.05」という値が記憶されている。また、図5(a)～(c)に示す例では、分野(カテゴリ)として、スポーツ、経済、政治、の3分野について示され、それぞれの学習用文書の文書数は、80文書、100文書、150文書であるものとする。更に、図5(a)～(c)に示す各カテゴリに属する語Wは、各カテゴリに属する文書の中に存在する語であり、その数値は、文書に含まれる語の頻度を示す。ここで、或る語の頻度として、例えば「その語のカテゴリに出現する個数の総計」や「カテゴリ内の、その語を含む文書の個数」など、任意の指標を採用することができる。

【0032】

図5(a)に示すように、まず、カテゴリ「スポーツ」の中で頻度が50である

単語「日本」を不要語とするか否かの判断を行なう。従来では、ただ単に、この頻度50が大きい小さいだけを判断対象としていたが、本実施の形態では、他のカテゴリにおける頻度の状況を確認し、カテゴリ間の相対的な出現頻度をもとに不要語を決定している。そのために、他のカテゴリである「経済」の文書の中で、単語「日本」がどの程度、使用され、出現しているか、を判断している。より具体的には、カテゴリ「経済」の文書数を閾値Rで掛け合わせた値($100 \times 0.05 = 5$)と、単語「日本」の頻度(30個)とが比較される。30は5よりも大きい($30 > 5$)ことから、「スポーツ」で用いられる単語「日本」は、他のカテゴリ(例えば「経済」)でも頻繁に用いられる可能性がある単語と考えられる。従って、実際に文書の分類作業を行なう際、「スポーツ」のカテゴリを判断する判断対象として「日本」は好ましくないと考えられる。そこで、カテゴリ「スポーツ」では、単語「日本」が不要語とされる。

【0033】

次に、図5(b)に示すように、カテゴリ「スポーツ」の中で単語「代表」を不要語とするか否かの判断を行なう。まず、他のカテゴリの1つである「経済」の中で、単語「代表」の頻度は2であり、カテゴリ「経済」の文書数を閾値Rで掛け合わせた値($100 \times 0.05 = 5$)と比較して小さい($2 < 5$)ことから、この段階では、カテゴリ「スポーツ」での不要語とは判断しない。しかしながら、もう一方の他のカテゴリである「政治」の中で、単語「代表」の頻度は8である。このとき、カテゴリ「政治」の文書数を閾値Rで掛け合わせた値($150 \times 0.05 = 7.5$)と比較して、出現頻度が大きい($8 > 7.5$)ことが理解できる。その結果、カテゴリ「スポーツ」の中の単語「代表」は、他のカテゴリの状態を判断して、識別単語としては好ましいものとは判断できない。そこで、「スポーツ」の中の単語「代表」は、不要語であると判断される。

【0034】

更に、図5(c)に示すように、カテゴリ「スポーツ」の中で単語「選手」を不要語とするか否かの判断を行なう。まず、他のカテゴリの1つであるカテゴリ「経済」の中で、単語「選手」の頻度は3であり、カテゴリ「経済」の文書数を閾値Rで掛け合わせた値($100 \times 0.05 = 5$)と比較して小さい($3 < 5$)ことか

ら、単語「選手」は、カテゴリ「スポーツ」での不要語とは判断しない。また、もう一方の他のカテゴリであるカテゴリ「政治」の中で、単語「選手」の頻度は1である。カテゴリ「政治」の文書数を閾値Rで掛け合わせた値($150 \times 0.05 = 7.5$)と比較して小さい($1 < 7.5$)ことが理解できる。従って、カテゴリ「スポーツ」の中の単語「選手」は、他のカテゴリにて出現頻度が低く、識別単語として好ましいものと判断され、「スポーツ」の中の単語「選手」は不要語ではなく、削除されずに残される。

【0035】

図6は、図5(a)～(c)によって、全てのカテゴリに対して不要語を除去した後の状態を説明するための図である。上述したようなアルゴリズムによる不要語除去作業を全てのカテゴリに対して施す。図6において、斜線で示される領域に存在する単語は、不要語として除去される単語である。カテゴリ「スポーツ」では単語「日本」、「代表」が、カテゴリ「経済」では単語「日本」、「選手」、「代表」が、カテゴリ「政治」では単語「日本」、「代表」、「銀行」、「選手」が、不要語として除去される。

【0036】

図7は、図3にて示した、分野テーブル作成部31にて作成されメモリに格納されるテーブル例から、不要語を除去した後の分野テーブルの例を示した図である。図3と同様に、「スポーツ」の分野を例に挙げている。テーブル情報は、不要語除去後に残った単語を特定する番号である単語IDごとに、単語、単語の品詞、単語の出現頻度が表されている。尚、図3と同様に、単語の出現頻度は、「学習用の文書集合中に出現した回数の総計」を表している。図7に示すような、不要語決定・除去部32にて不要語が除去された分野テーブルは、分類用カタログ格納装置22に分類用カタログとして格納される。尚、分類用カタログ格納装置22に格納するに際し、図7に示すような、不要語が取り除かれた語のリストをそのまま用いてもよく、または、このリストに既存の「語の重み付け手法」を用いてリストを改良して格納することもできる。

【0037】

以上のようにして不要語が除去された結果を用いて、実際に文書分類処理が実

行される。不要語を除去して得られた分野テーブルを文書分類処理に適用する方法としては幾つかのものが考えられるが、ここでは、「ベクトル空間モデル」と呼ばれる方法を例に挙げて説明する。

【0038】

分類用カタログ格納装置22には、不要語除去を経て作成された分野テーブルが格納されているが、分野(カテゴリ)ごとに、「語」と「語の重み」のペアが登録されている。図6に示す例では、「スポーツ」のカテゴリにて、語である「選手」、語の重みとして「20」が登録されている。例えば、図6に示すような場合には、「選手」、「取引」、「銀行」、「ビール」、「首相」という5個の単語(語)の組を基底とするベクトル空間を考え、この空間内で「文書と各分野との距離」を計算する。尚、複数の分野に出てくる場合には、重複して出てきている語をまとめて1個としてベクトル空間が作成される。図6に示す例では、各分野のベクトルは、以下のようになる。

スポーツ : (20, 0, 0, 0, 0)

経済 : (0, 20, 10, 3, 0)

政治 : (0, 0, 0, 0, 100)

【0039】

次に、分類対象の文書から、文書ベクトルを作成する方法を説明する。ここでは、まず、分類対象文書格納装置23から得られる分類対象の文書Dを形態素解析し、単語とその出現頻度との組をテーブルにする。例えば、

分類対象の文書の内容：

「A国の首相が、B国の首相とイラク問題について会談した。」

について、形態素解析を行い、下記のようなテーブルを作成する。

(A, 1)、(国, 2)、(首相, 2)、(イラク, 1)、(問題, 1)、(会談, 1)

次に、このようにして作成されたテーブルと、既に作成されているベクトル空間の基底とを比較し、ベクトル空間の基底になっている(登録されている)単語の情報のみを用いてベクトルを作成し、分類対象文書のベクトルが生成される。上記の例では、生成される文書ベクトルは、

選手、 取引、 銀行、 ビール、 首相

(0、 0、 0、 0、 2)

となる。

【0040】

その後、「文書と各分野との距離」の計算に、以上のようにして生成されたベクトルのなす角度の余弦が用いられる。

図8(a),(b)は、本実施の形態にて利用されるベクトル空間モデルを説明するための図である。この余弦は、図8(a)に示すベクトルAとベクトルBに対して、その角度を θ とすると、

$$\cos \theta = (A \cdot B) \div (|A| |B|)$$

で定義される。ここで、 $A \cdot B$ はAとBとの積、 $|A|$ はAのノルム(長さ)を表す。余弦の値、即ち $\cos \theta$ は、0と1の間をとり、1に近いほど θ が小さくなる。つまり、 $\cos \theta$ の値が大きいほど、AとBとは互いに「近い」と考えられる。

【0041】

文書の分類において、余弦は、次のようにして用いることができる。分類したい文書に対応するベクトルをA、分野に対応するベクトルをBとし、各Bに対して、AとBとの余弦を計算する。Aに対して余弦の値を最も大きくするようなBの分野を、Aが属する分野と判定すればよい。図8(b)に示すように、分類対象文書をベクトルAにとり、政治、経済、スポーツの各カテゴリをベクトルBにとる。そして、分類対象文書と政治、分類対象文書と経済、分類対象文書とスポーツ、の各々の余弦を、上述した式にて算出する。図8(b)に示す例では、分類対象文書と政治との角度が最も小さく、余弦が最も大きくなり、分類対象文書は「政治」のカテゴリに属するものと判定することができる。

【0042】

図9は、このようなベクトル空間モデルを用いて文書分類処理部33にて実行される文書分類処理の流れを示したフローチャートである。文書分類処理部33では、まず、分類対象文書格納装置23から分類対象文書Dが取得される(ステップ301)。次に、分類対象文書Dの単語を全て抽出し、分類対象文書Dに対応するベクトルV_dを作成する(ステップ302)。ここで、全ての分野について

作業したか否かが判断され(ステップ303)、作業が残っている場合には、分野を1つ選んでAとする(ステップ304)。そして、ベクトルV_dと、Aに対応するベクトルV_aとの距離を、上述のようにして計算する(ステップ305)。ステップ303へ戻り、全ての作業について終了した場合には、計算した距離を用いて、分類対象文書Dの分類先の分野を決定し(ステップ306)、分類結果格納装置24に結果を格納して処理が終了する。

【0043】

以上、詳述したように、本実施の形態では、文書自動分類における不要語を「他のどれかのカテゴリにもある程度以上含まれている語」という定義を行い、カテゴリ間の相対的な出現頻度から不要語除去を行なっている。これによって、カテゴリの特徴付けに役立たない語(不要語)を新たに定義することができ、この定義によって、従来の手法に比べて、より効果的に不要語を除去することができる。また、不要語が除去されたリストを分類用カタログ格納装置22に格納し、このリストを用いて実際の文書分類処理を実行することで、実際の文書処理に際して不要語か否かを判断するといった手間を省くことができる。即ち、実際の分類対象文書を解析して不要語を除去する必要がなく、分類作業を迅速化することが可能となる。

【0044】

【発明の効果】

以上説明したように、本発明によれば、文書の自動分類において、不要語を効果的に取り除くことが可能となる。

【図面の簡単な説明】

【図1】 本実施の形態が適用される文書自動分類システムの構成を示したブロック図である。

【図2】 分野テーブル作成部にてなされる処理を示したフローチャートである。

【図3】 図2にて説明したような分野テーブル作成部にて作成されメモリに格納されるテーブル例を示した図である。

【図4】 不要語除去部にてなされる処理を示したフローチャートである。

【図 5】 (a)～(c)は、この不要語処理のアルゴリズムについて更に詳しく説明するための図である。

【図 6】 図 5(a)～(c)によって、全てのカテゴリに対して不要語を除去した後の状態を説明するための図である。

【図 7】 図 3 にて示した分野テーブル作成部にて作成されメモリに格納されるテーブル例から不要語を除去した後の分野テーブルの例を示した図である。

【図 8】 (a),(b)は、本実施の形態にて利用されるベクトル空間モデルを説明するための図である。

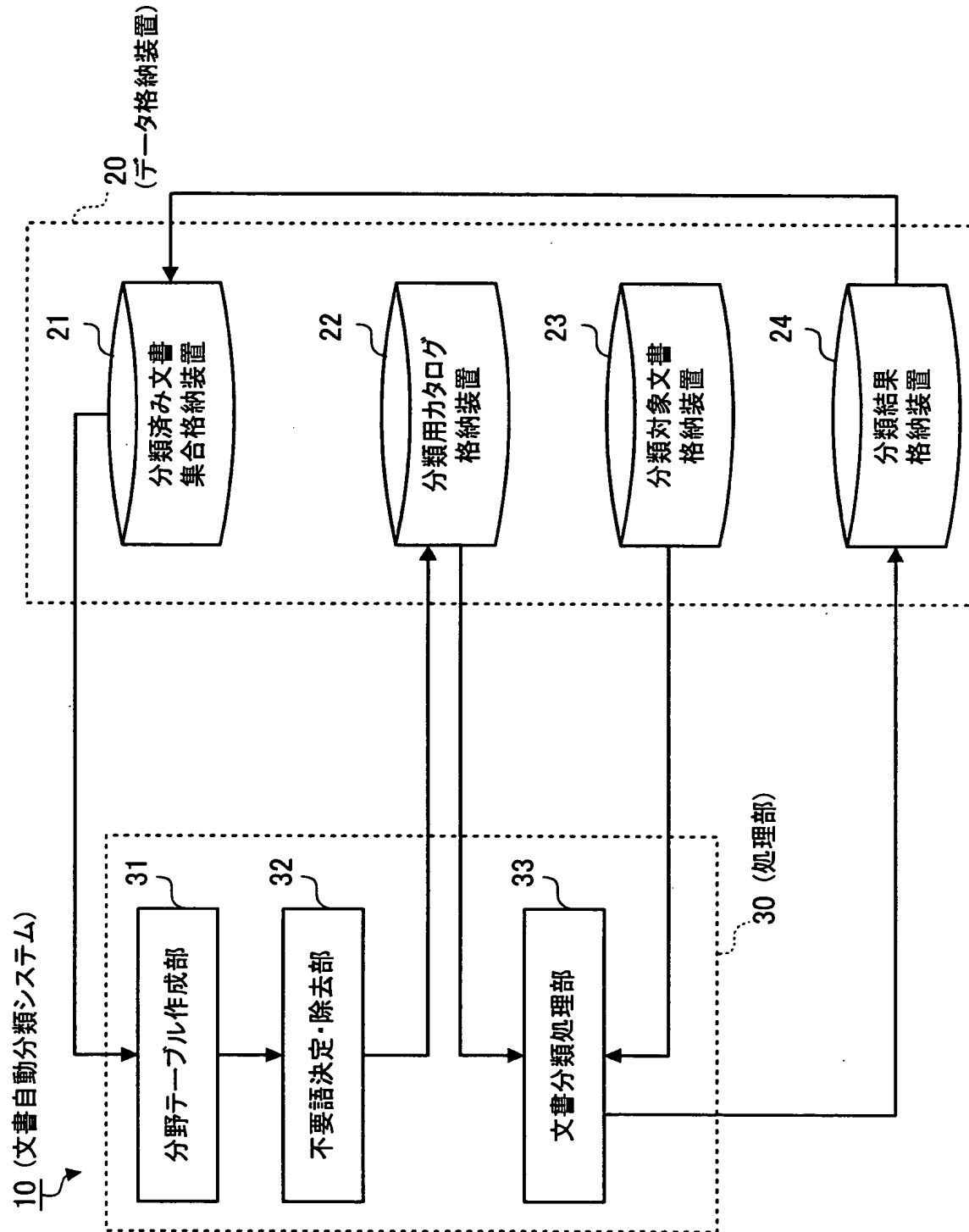
【図 9】 ベクトル空間モデルを用いて文書分類処理部にて実行される文書分類処理の流れを示したフローチャートである。

【符号の説明】

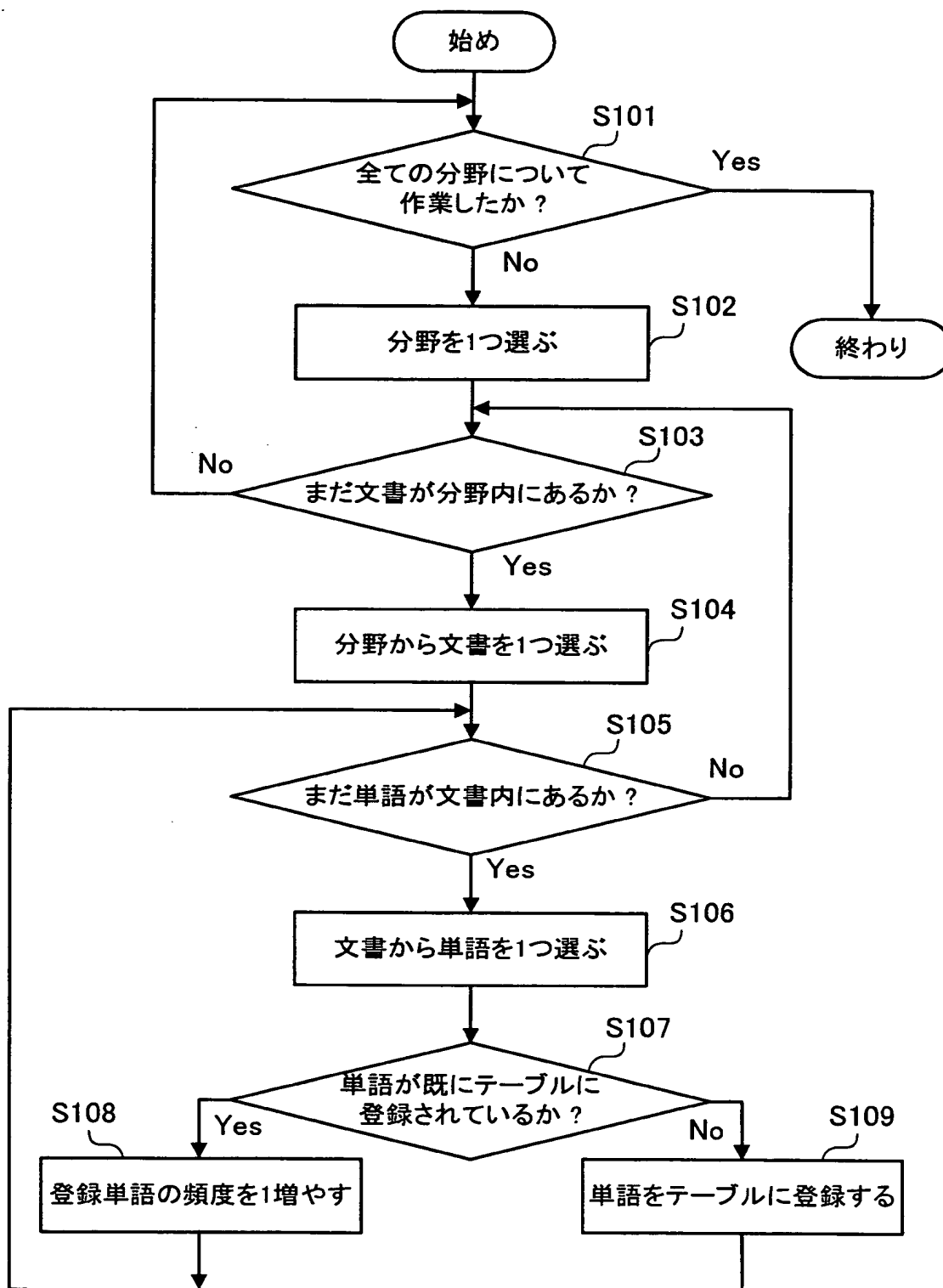
1 0 …文書自動分類システム、2 0 …データ格納装置、2 1 …分類済み文書集合格納装置、2 2 …分類用カタログ格納装置、2 3 …分類対象文書格納装置、2 4 …分類結果格納装置、3 0 …処理部、3 1 …分野テーブル作成部、3 2 …不要語決定・除去部、3 3 …文書分類処理部

【書類名】 図面

【図 1】



【図 2】



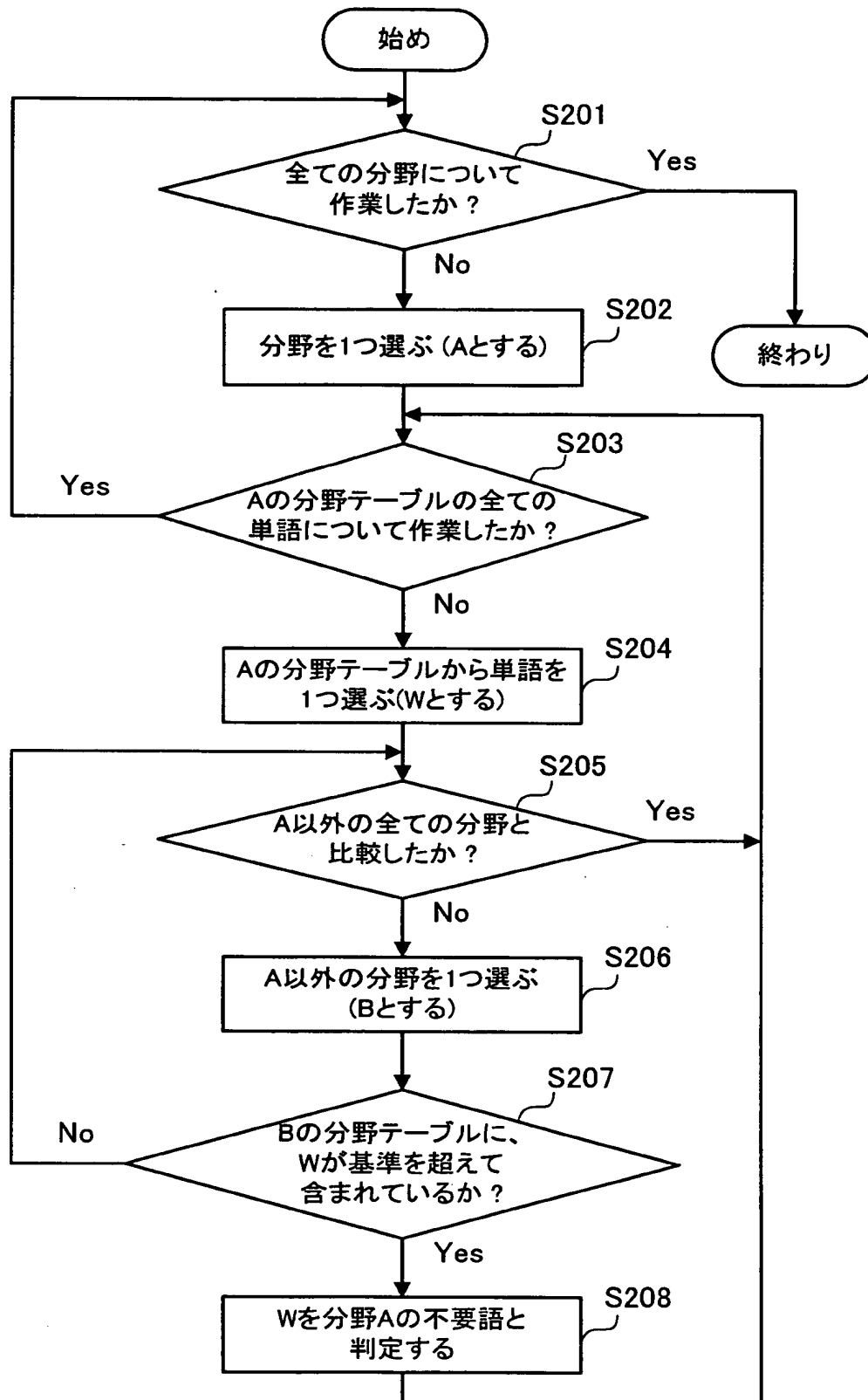
【図3】

分野テーブル(不要語除去前)

単語ID 品詞 出現頻度
 ↓ ↓ ↓
 単語 単語 単語

0: れ: 動詞: 8349	26: 世界: 名詞: 252
1: い: 動詞: 4841	27: 話: 動詞: 238
2: な: 動詞: 4335	28: 敗れ: 動詞: 234
3: 行: 動詞: 3570	29: 女子: 名詞: 221
4: 日: 名詞: 2250	30: 投手: 名詞: 220
5: 日本: 固有名詞: 1856	31: ついて: 動詞: 220
6: し: 動詞: 1716	32: 監督: 名詞: 216
7: チーム: 名詞: 1175	33: 外野手: 名詞: 216
8: 選手: 名詞: 1128	34: ○○○○: 固有名詞: 216
9: 試合: 名詞: 714	35: よ: 動詞: 208
10: い: 動詞: 660	36: 安打: 名詞: 204
11: 大: 名詞: 583	37: 大リーグ: 名詞: 195
12: 通算: 名詞: 532	38: 世界選手権: 名詞: 192
13: 決勝: 名詞: 513	39: 東京: 固有名詞: 190
14: 優勝: 名詞: 494	40: 選手権: 名詞: 182
15: 出場: 動詞: 456	41: 出場: 名詞: 182
16: リーグ: 名詞: 448	42: シドニー五輪: 固有名詞: 176
17: 大会: 名詞: 432	43: 首位: 名詞: 168
18: 打席: 名詞: 396	45: せ: 動詞: 154
19: あ: 動詞: 357	46: ○○○: 固有名詞: 153
20: 連続: 名詞: 357	:
21: 男子: 名詞: 345	:
22: 決め: 動詞: 323	:
23: 代表: 名詞: 312	(以下、同様にして続く)
24: 打: 動詞: 306	
25: 今季: 名詞: 266	

【図 4】



【図 5】

(a) スポーツの「日本」を不要語とするか否かの判断

閾値 $R = 0.05$

$80 \times 0.05 = 4$		$100 \times 0.05 = 5$		$150 \times 0.05 = 7.5$	
スポーツ (80文書)		経済 (100文書)		政治 (150文書)	
日本	50	日本	30	首相	100
代表	20	取引	20	日本	70
選手	20	銀行	10	代表	8
		選手	3	銀行	5
		ビール	3	選手	1
		代表	2		

30 > 5

(b) スポーツの「代表」を不要語とするか否かの判断

閾値 $R = 0.05$

$80 \times 0.05 = 4$		$100 \times 0.05 = 5$		$150 \times 0.05 = 7.5$	
スポーツ (80文書)		経済 (100文書)		政治 (150文書)	
日本	50	日本	30	首相	100
代表	20	取引	20	日本	70
選手	20	銀行	10	代表	8
		選手	3	銀行	5
		ビール	3	選手	1
		代表	2		

2 < 5

8 > 7.5

(c) スポーツの「選手」を不要語とするか否かの判断

閾値 $R = 0.05$

$80 \times 0.05 = 4$		$100 \times 0.05 = 5$		$150 \times 0.05 = 7.5$	
スポーツ (80文書)		経済 (100文書)		政治 (150文書)	
日本	50	日本	30	首相	100
代表	20	取引	20	日本	70
選手	20	銀行	10	代表	8
		選手	3	銀行	5
		ビール	3	選手	1
		代表	2		

3 < 5

1 < 7.5

【図 6】

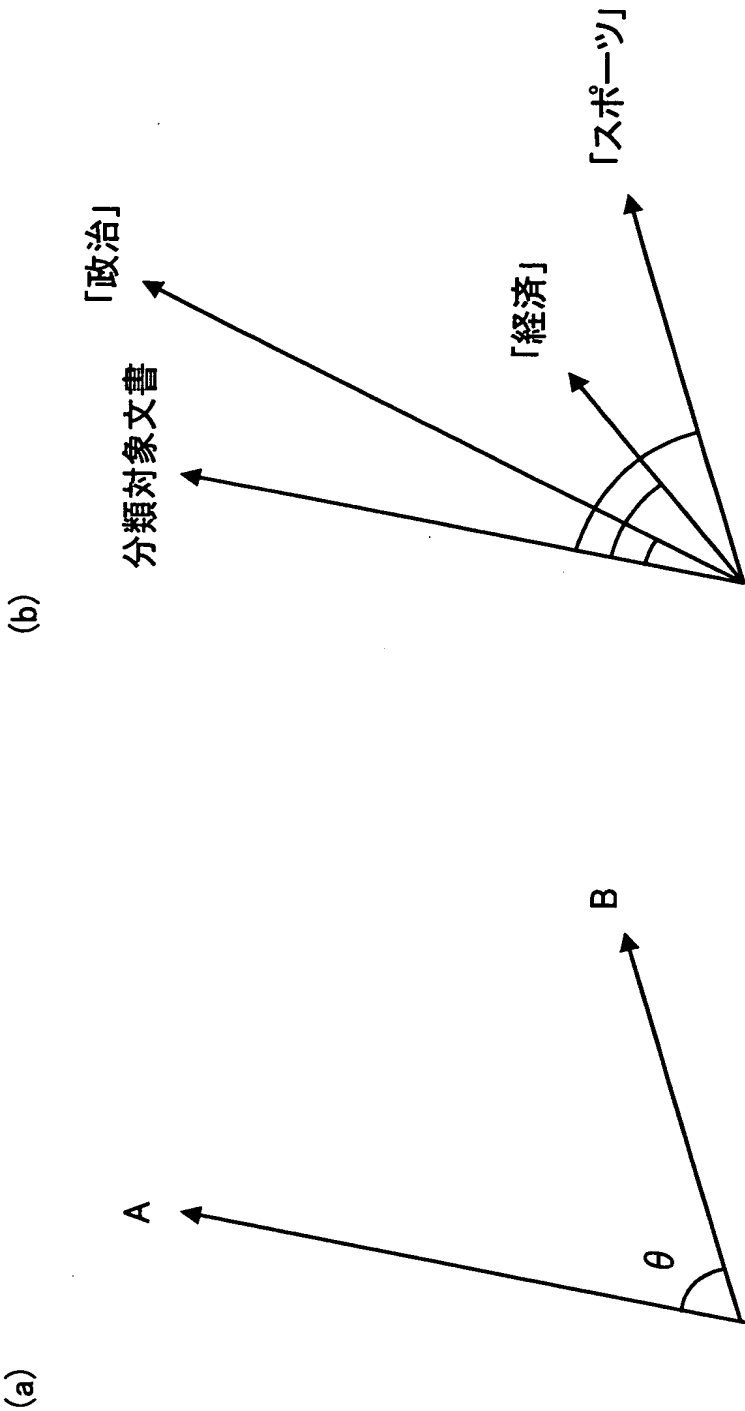
スポーツ (80文書)	経済 (100文書)	政治 (150文書)
日本 50	日本 30	首相 100
代表 20	取引 20	日本 10
選手 20	銀行 10	代表 8
	選手 3	銀行 5
	ビール 3	選手 1
	代表 2	

【図 7】

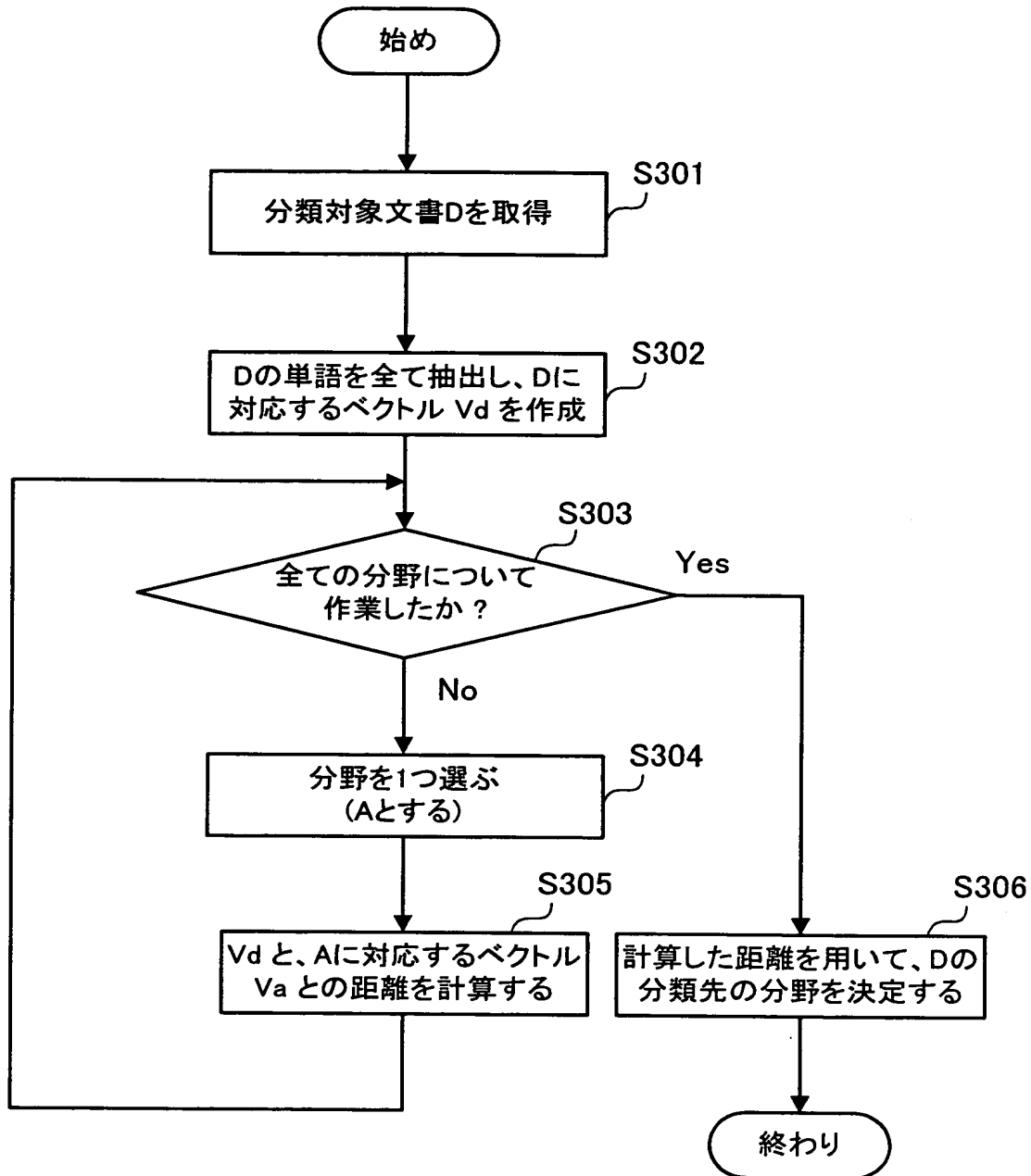
分野テーブル(不要語除去後)

単語ID	単語	品詞	出現頻度
0	選手	名詞	1128
1	試合	名詞	714
2	決勝	名詞	513
3	出場	動詞	456
4	リーグ	名詞	448
5	打席	名詞	396
6	男子	名詞	345
7	今季	名詞	266
8	敗れ	動詞	234
9	女子	名詞	221
10	投手	名詞	220
11	外野手	名詞	216
12	〇〇〇〇	固有名詞	216
13	安打	名詞	204
14	大リーグ	名詞	195
15	世界選手権	名詞	192
16	出場	名詞	182
17	選手権	名詞	182
18	シドニー五輪	type 145	176
19	首位	名詞	168
20	〇〇〇	固有名詞	153
:			
:			
:			
(以下同様に続く)			

【図 8】



【図9】



【書類名】 要約書

【要約】

【課題】 文書自動分類において、不要語を効果的に取り除く。

【解決手段】 分野ごとに分類済みの文書を格納する分類済み文書集合格納装置 2 1 と、この分類済み文書集合格納装置 2 1 から取得された文書に含まれる単語の出現頻度の情報を含む分野別のテーブルを作成する分野テーブル作成部 3 1 と、この分野テーブル作成部 3 1 により作成された分野別のテーブルから得られる所定の語における各分野での出現頻度に基づいて、テーブルから分野別に不要語を除去する不要語決定・除去部 3 2 と、この不要語決定・除去部 3 2 により不要語が除去されたテーブルを格納する分類用カタログ格納装置 2 2 と、分類される分類対象文書を格納する分類対象文書格納装置 2 3 と、この分類対象文書格納装置 2 3 に格納された分類対象文書に対し、分類用カタログ格納装置 2 2 に格納されたテーブルを用いて分類処理を行なう文書分類処理部 3 3 とを含む。

【選択図】 図 1

認定・付加情報

特許出願の番号	特願 2002-301539
受付番号	50201554965
書類名	特許願
担当官	土井 恵子 4264
作成日	平成14年10月17日

<認定情報・付加情報>

【特許出願人】

【識別番号】	390009531
【住所又は居所】	アメリカ合衆国10504、ニューヨーク州 アーモンク ニュー オーチャード ロード
【氏名又は名称】	インターナショナル・ビジネス・マシーンズ・コーポレーション

【代理人】

【識別番号】	100086243
【住所又は居所】	神奈川県大和市下鶴間1623番地14 日本アイ・ビー・エム株式会社 大和事業所内
【氏名又は名称】	坂口 博

【代理人】

【識別番号】	100091568
【住所又は居所】	神奈川県大和市下鶴間1623番地14 日本アイ・ビー・エム株式会社 大和事業所内
【氏名又は名称】	市位 嘉宏

【代理人】

【識別番号】	100108501
【住所又は居所】	神奈川県大和市下鶴間1623番14 日本アイ・ビー・エム株式会社 知的所有権
【氏名又は名称】	上野 剛史

【復代理人】

申請人	
【識別番号】	100104880
【住所又は居所】	東京都港区赤坂5-4-11 山口建設第2ビル 6F セリオ国際特許事務所
【氏名又は名称】	古部 次郎

次頁無

出 願 人 履 歴 情 報

識別番号 [390009531]

1. 変更年月日 2002年 6月 3日

[変更理由] 住所変更

住 所 アメリカ合衆国10504、ニューヨーク州 アーモンク ニ
ュー オーチャード ロード

氏 名 インターナショナル・ビジネス・マシーンズ・コーポレーショ
ン